# Intelligent Automation Incorporated

# Enhancements for a Dynamic Data Warehousing and Mining System for Large-Scale Human Social Cultural Behavioral (HSBC) Data

## Final Report

Reporting Period: March 22, 2016 – September 18, 2016

Contract No.  N00014-16-P-3014

Prepared by

Onur Savas, Ph.D.

# REPORT DOCUMENTATION PAGE

The public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing the burden, to Department of Defense, Washington Headquarters Services, Directorate for Information Operations and Reports (0704-0188), 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to any penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.
**PLEASE DO NOT RETURN YOUR FORM TO THE ABOVE ADDRESS.**

| 1. REPORT DATE *(DD-MM-YYYY)* | 2. REPORT TYPE | 3. DATES COVERED *(From - To)* |
|---|---|---|
| 26/09/2016 | Final Report | March 22, 2016 – September 18, 2016 |

**4. TITLE AND SUBTITLE**

Enhancements for a Dynamic Data Warehousing and Mining System for Large-Scale Human Social Cultural Behavioral (HSBC) Data

**5a. CONTRACT NUMBER**

N00014-16-P-3014

**5b. GRANT NUMBER**

**5c. PROGRAM ELEMENT NUMBER**

**6. AUTHOR(S)**

Onur Savas

**5d. PROJECT NUMBER**

**5e. TASK NUMBER**

**5f. WORK UNIT NUMBER**

**7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES)**

Intelligent Automation, Inc.
15400 Calhoun Drive, Suite 190
Rockville, MD 20855

**8. PERFORMING ORGANIZATION REPORT NUMBER**

**9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)**

Office of Naval Research (BD253)
875 N. Randolph Street
Arlington, VA 22203-1995

**10. SPONSOR/MONITOR'S ACRONYM(S)**

**11. SPONSOR/MONITOR'S REPORT NUMBER(S)**

**12. DISTRIBUTION/AVAILABILITY STATEMENT**

DISTRIBUTION A: Approved for public release; distribution is unlimited.

**13. SUPPLEMENTARY NOTES**

Report contains color

**14. ABSTRACT**

In this project, to continue our development of a system that can dynamically collect and warehouse unfiltered textual communication data and make this data available to support HSCB modeling, decision making, and course of action development, we have (i) developed algorithms and an end-to-end computational framework to support graph analysis and visualization, (ii) developed graph querying capabilities such as top K and n-hop neighborhood, (iii) developed automated YouTube video metadata collection and analysis of YouTube data, and (iv) developed automated VK data collection and analysis of VK data. The features are released as part of Scraawl 2.0.

**15. SUBJECT TERMS**

social media, HSCB, data warehousing, large-scale analysis, OSINT

| 16. SECURITY CLASSIFICATION OF: | | | 17. LIMITATION OF ABSTRACT | 18. NUMBER OF PAGES | 19a. NAME OF RESPONSIBLE PERSON |
|---|---|---|---|---|---|
| a. REPORT | b. ABSTRACT | c. THIS PAGE | SAR | 15 | Dr. Rebecca Goolsby |
| Unclassified | Unclassified | Unclassified | | | 19b. TELEPHONE NUMBER *(Include area code)* (703) 588-0558 |

# Table of Contents

# 1   Executive Summary

In this project, to continue our development of a system that can dynamically collect and warehouse unfiltered textual communication data and make this data available to support HSCB modeling, decision making, and course of action development, we have accomplished the following.

**First,** we have developed algorithms and an end-to-end computational framework to ingest Twitter data, convert it into Twitter interaction and/or friendship graphs, and store in a graph database, namely OrientDB. We have also enhanced Scraawl's current graph visualization for scalability and usability.

**Second,** we have developed graph querying capabilities and implemented various UIs to visualize the outputs of these queries. In particular, we have developed top K user and/or hashtag query, and n-hop neighborhood queries. Various attributes of the nodes and edges can also be interactively visualized.

**Third,** we have developed automated YouTube video metadata collection capabilities. In particular, we have developed two APIs for continuous (Streaming) and recent YouTube data searches, and designed a user-friendly UI for these searches using YouTube Data API v3. Algorithms to perform basic statistics computation, NER, and map visualization have also been designed and implemented.

**Fourth,** we have developed automated VK data collection capabilities. Algorithms to perform basic statistics computation, NER, and map visualization have also been designed and implemented.

**Finally,** we released Scraawl 2.0, which incorporates these features.

## 2 Technical Work

### 2.1 Store Twitter Interaction Graphs in OrientDB

The first step in developing the querying capability is to store a graph efficiently. A graph database, for all practical purposes, can be represented as the graph itself. In its simplest form, we consider the *Twitter interaction graph* modeled as an undirected graph $G = (V, E)$, where $V$ is the set of nodes (vertices) and $E$ is the set of edges. For a collection of tweets $\{\tau(\theta)| \theta \in \mathbb{Z}^+\}$, where each tweet $\tau(.)$ can be uniquely identified by its unique *tweet ID* $\theta \in \mathbb{Z}^+$, let $\tau(\theta)$ be tweeted by user $u_{\tau(\theta)}$ and let $u_{\tau(\theta)}$ have retweeted, mentioned, or replied to $K_\theta$ users $\mathcal{I}(\tau(\theta)) = \{v^1_{\tau(\theta)}, v^2_{\tau(\theta)}, ..., v^{K_\theta}_{\tau(\theta)}\}$. Of course, if no retweets, mentions, or replies are present, then $\mathcal{I}(\tau(\theta)) = \emptyset$. We can then unambiguously specify the Twitter Interaction Graph $G$ with $V = \{u_{\tau(\theta)} \cup \mathcal{I}(\tau(\theta))| \theta \in \mathbb{Z}^+\}$ and $E = \{u_{\tau(\theta)} \times \mathcal{I}(\tau(\theta))| \theta \in \mathbb{Z}^+\} = \{(u_{\tau(\theta)}, v^1_{\tau(\theta)}), (u_{\tau(\theta)}, v^2_{\tau(\theta)}), ..., (u_{\tau(\theta)}, v^K_{\tau(\theta)})| \theta \in \mathbb{Z}^+\}$. One can enhance $G$ by adding other types of interactions, e.g., by adding an edge if a user uses a hashtag in his/her tweet. In this particular case, the vertex set will have hashtags as well. Note that this definition can trivially be extended for *friendship graphs*, where the interactions between users/nodes are whether they are friends and/or followers.

To implement the graph modeling capabilities above, we used an open source graph database, namely OrientDB (http://www.orientdb.com), and an SQL-like graph querying language. OrientDB provides an NoSQL engine that stores and queries graphs via both (i) a graph database API and document API, and (ii) supports schema-less, schema-full, and schema-mixed modes. Our initial experimentation with OrientDB, an open source graph database, and its SQL-like graph querying language yielded promising results. In particular, we inserted a Twitter Interaction Graph with $|V| > 34k$ and $|E| > 421k$ in less than 6 minutes. In addition, we ordered nodes of this graph by degree in less than 2 seconds. A basic graph service that performs basic graph database management operations, e.g., insert, delete, update, is also implemented.

We have also designed and implemented the *k*-hop neighbor queries. In brief and informally, starting from a node *v*, *k*-hop queries finds the neighbors of *v*, neighbors of neighbors of *v* (a.k.a. second degree neighbors), and so on until *k*-hop neighbors. In particular, for a querying function *Q(.)*, a "2-hop" neighbors query is implemented as follows.

$$Q(v) = \{(v, v') \cup (v', v'')| (v, v') \in E \text{ and } (v', v'') \in E \text{ for } \forall v', v'' \in V\}.$$

This can be generalized to *k*-hop by including neighbors of neighbors of neighbors up *k*-hop.

This querying service has also been implemented using a Groovy/Grails framework with functions written in Java.

## 2.2 Design and Implement Visualization and UI

To visualize the graph returned by the queries, we have further developed IAI's graph visualization capabilities. We have also incorporated a UI to call the *k*-hop query service. These capabilities are all incorporated to Scraawl.

Figure 1 is a visualization for 4-hop neighborhood query starting from #EasterEggRoll. The Twitter interaction graph has been created from a collection of tweets that have been collected between 28 Mar, 2016 03:00 AM and 29 Mar, 2016 02:59 using the keyword #EasterEggRoll. Overall, 25063 tweets were collected. The graph has 16171 nodes (vertices) and 48594 edges.

The "Neighbors" button to the upper left corner call the k-hop querying service with (i) starting node (in this case #EasterEggRoll), and (ii) the *k*-parameter (in this case *k* = 4).



**Figure 1: 4-hop neighborhood query starting from "#eastereggroll".**

## 2.3 Top K User and Hashtag Subgraph Querying and Visualization

### 2.3.1 Top K Query Implementation

As reported in Progress Report No. 1, we consider the Twitter interaction graph modeled as an undirected graph $G = (V, E)$, where $V$ is the set of nodes (vertices) and $E$ is the set of edges. For a collection of tweets $\{\tau(\theta) | \theta \in \mathbb{Z}^+\}$, where each tweet $\tau(.)$ can be uniquely identified by its unique *tweet ID* $\theta \in \mathbb{Z}^+$, let $\tau(\theta)$ be tweeted by user $u_{\tau(\theta)}$ and let $u_{\tau(\theta)}$ have retweeted, mentioned, or replied to $K_\theta$ users $\mathcal{I}(\tau(\theta)) = \{v^1_{\tau(\theta)}, v^2_{\tau(\theta)}, \dots, v^{K_\theta}_{\tau(\theta)}\}$. Of course, if no retweets, mentions, or replies are present, then $\mathcal{I}(\tau(\theta)) = \emptyset$. We can then unambiguously specify the Twitter Interaction Graph $G$ with $V = \{u_{\tau(\theta)} \cup \mathcal{I}(\tau(\theta)) | \theta \in \mathbb{Z}^+\}$ and $E = \{u_{\tau(\theta)} \times \mathcal{I}(\tau(\theta)) | \theta \in \mathbb{Z}^+\} = \{(u_{\tau(\theta)}, v^1_{\tau(\theta)}), (u_{\tau(\theta)}, v^2_{\tau(\theta)}), \dots, (u_{\tau(\theta)}, v^K_{\tau(\theta)}) | \theta \in \mathbb{Z}^+\}$.

Our interest lies in the querying of this graph. We first define a generic query operator over the graph that will return a subset of the graph, i.e., $Q(.): G \mapsto G$. In this reporting period, we have designed and implemented top[1] K neighborhood queries. Formally, we define the top K neighborhood query as

$$Q\left(\bigcup_i v_i\right) = G(V', E') \; where \; V' = \left(\bigcup_i v_i \cup N(v_i)\right) and \; s.t. \; E' = \bigcup_{i,j}(v_i, v_j(i)),$$
$$where \; v_j(i) \in N(v_i)$$

for a set of "top" $\bigcup_i v_i$ vertices and $N(v_i)$ denotes the immediate neighborhood of $v_i$. Recall from Report 1 that we have designed and implemented a graph querying system using the graph database OrientDB. We have added the above querying capability implemented using OrientDB and made it available through Scraawl.

In particular, we have implemented top K neighborhood queries for the following top K.

**1. Top K Users:** The top tweeting users.

**2. Top K Hashtags:** The top tweeted hashtags.

**3. Top Connected Users:** Top users ranked by degree in the Social Graph.

**4. Top Connected Hashtags**: Top hashtags ranked by degree in the Social Graph.

We have also enhanced the visualization for top K queries. We have improved the visualization by (i) automatically adjusting zooming and translating so that the graph fits into the visible screen, (ii) incorporating web workers, i.e., background computation threads, so that graph layout is computed without interfering with user's interaction with the page, and (iii) adjusted the parameters of the graph layout so that there is less

---

[1] As defined according to context as shown in Section 1.1.2.

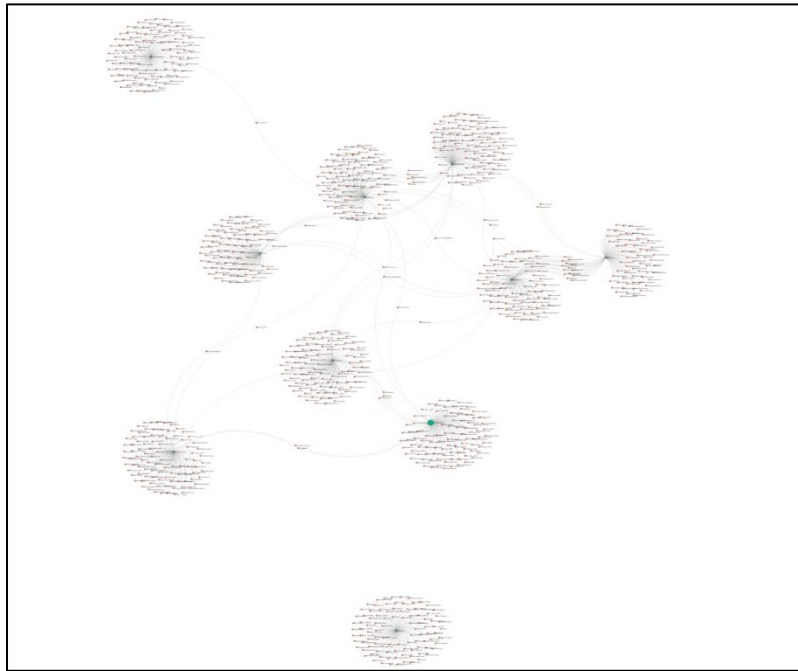cluttering. Figure 2 is a representative visualization for a top 10 query as integrated into Scraawl.



**Figure 2: Representative visualization of Top 10 users.**

### 2.3.2   Graph Visualization for Friendship Graph

We have also designed and implemented a service that retrieves the friends and followers of a given user and visualizes it in Scraawl. Figure 3 is a representative graph visualization of the friendship graph of two users. The "queried users" are depicted as orange, the common friends/followers are shown as green, and the rest of the users are drawn in gray. Unlike the Social graph, the friendship graph is directed hence the direction of links are represented in the visualization. The relationship "being a friend" and "being a follower" are distinguished by using orange and gray colors.

**Figure 3: Representative Friendship Graph.**

## 2.4   YouTube Data Collection and Analytics

### 2.4.1   YouTube Data Collection
The first API allows for streaming searches, i.e., adds a query for continuous YouTube post collection. The second API allows searching on recent YouTube videos, and the resulting data will be saved in the configured database. Both APIs make use of the (i) `Search:list` Data API functionality, which returns a collection of search results that match the query parameters specified in the API request, and (ii) `Videos: list` Data API functionality, which returns a list of videos that match the API request parameters. Some of the major parameters used in the Scraawl streaming and recent API searches is shown in Table 1.

**Table 1: Scraawl API Major YouTube Data Collection Parameters.**

| Parameter | Explanation |
|---|---|
| **query.q** | The `q` parameter specifies the query term to search for. Two words separated by spaces can be treated as AND. Your request can also use the Boolean NOT (-) and OR (\|) operators to exclude videos or to find videos that are associated with one of several search terms. |
| **query.channelId** | The `channelId` parameter indicates that the API response |

| | should only contain resources created by the channel. |
|---|---|
| **query.location & query.locationRadius** | The `location` parameter, in conjunction with the `locationRadius` parameter, defines a circular geographic area and also restricts a search to videos that specify, in their metadata, a geographic location that falls within that area. |
| **query.publishedBefore** | The `publishedBefore` parameter indicates that the API response should only contain resources created before the specified time. |
| **query.publishedAfter** | The `publishedBefore` parameter indicates that the API response should only contain resources created before the specified time. |

We have also developed a UI to use the above APIs seamlessly. A representative UI is shown in Figure 4. The UI has the same look and feel with other social media searches, and can be accessed from "Create New Report" view under Scraawl. The UI allows to specify keywords with each box "AND"ed, and the translation capability using Google translate is integrated. The user can also choose between "Streaming" and "Recent" searches, which in turn calls one of the above APIs explained in Section 1.1.1. When "Recent" search is selected, the user has the option to select a time range. When "Streaming" search is selected, data collection will continue until a pre-specified time-out or the data collection limit is reached. Similar to other data feeds, we also allow the user to draw circular bubbles on the world map to restrict their searches to certain region(s) under "Additional Search Options."



**Figure 4: UI for YouTube Searches.**

### 2.4.2 YouTube Analytics

We have developed a fast and reliable computation of Top K statistics for YouTube videos. In particular, we compute Top Videos, Top Users, Top Words, and Top Languages. In all cases, "Top" refers to higher count. Similar to other data sources, we have also implemented a UI to show and interact with these statistics. Figure 5 shows a representative view of the UI. Top Videos, Top Users, and Top Words are shown as lists while Top Languages are shown as a pie chart. Every Top K Statistics also include a drill-down menu, where users can access using the "Details" button. In the drill-down menu, users can see Top 50 statistics and other relevant metrics. Also, users can filter to include or exclude the relevant posts.



**Figure 5: Representative Top K Statistics for YouTube Videos.**

Similar to the other data feeds, we have developed a timeline view that shows the number of posts. The timeline view is interactive and the selected portion can be filtered using "Filter By Range" as shown in Figure 6.



**Figure 6: Representative Timeline View.**

We show the snapshots of the Top 50 videos in a grid under "Media Gallery" view. The video snapshots are clickable and each click directs the user to a page that has statistics about video. A representative media gallery view is shown in Figure 7.
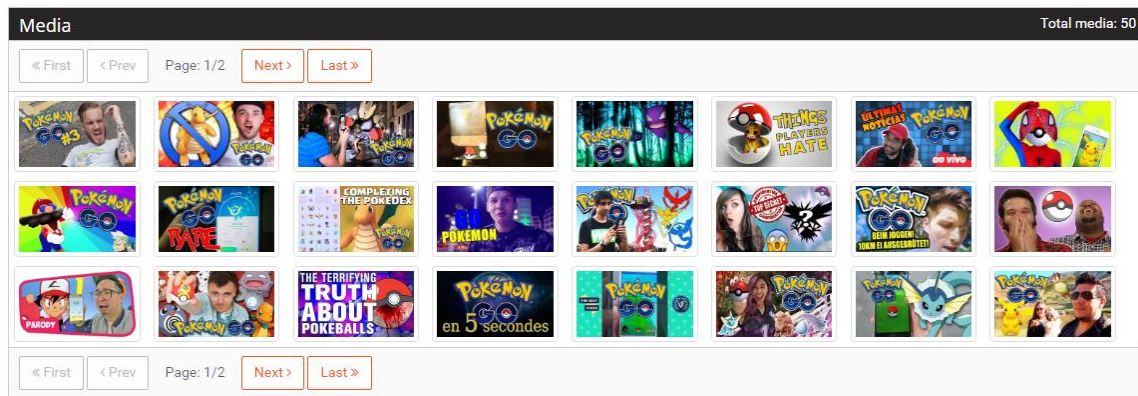
**Figure 7: Representative Media Gallery View.**

We perform Scraawl's NER algorithm to the text associated with YouTube post, which classifies the named entities into organizations, locations, and persons. Scraawl's NER also performs abbreviation extensions, e.g., UN is mapped to "United Nations." A representative NER view along with other statistics is shown in Figure 8.



**Figure 8: Representative NER View.**

We use Scraawl's Location Map and Heat Map view to display geolocation information about YouTube videos. In particular, we show geo-coded and geo-referenced YouTube posts. Geo-coded posts are those which have GPS information on them while geo-referenced posts are those which mention location information. A representative Map view is shown in Figure 9.
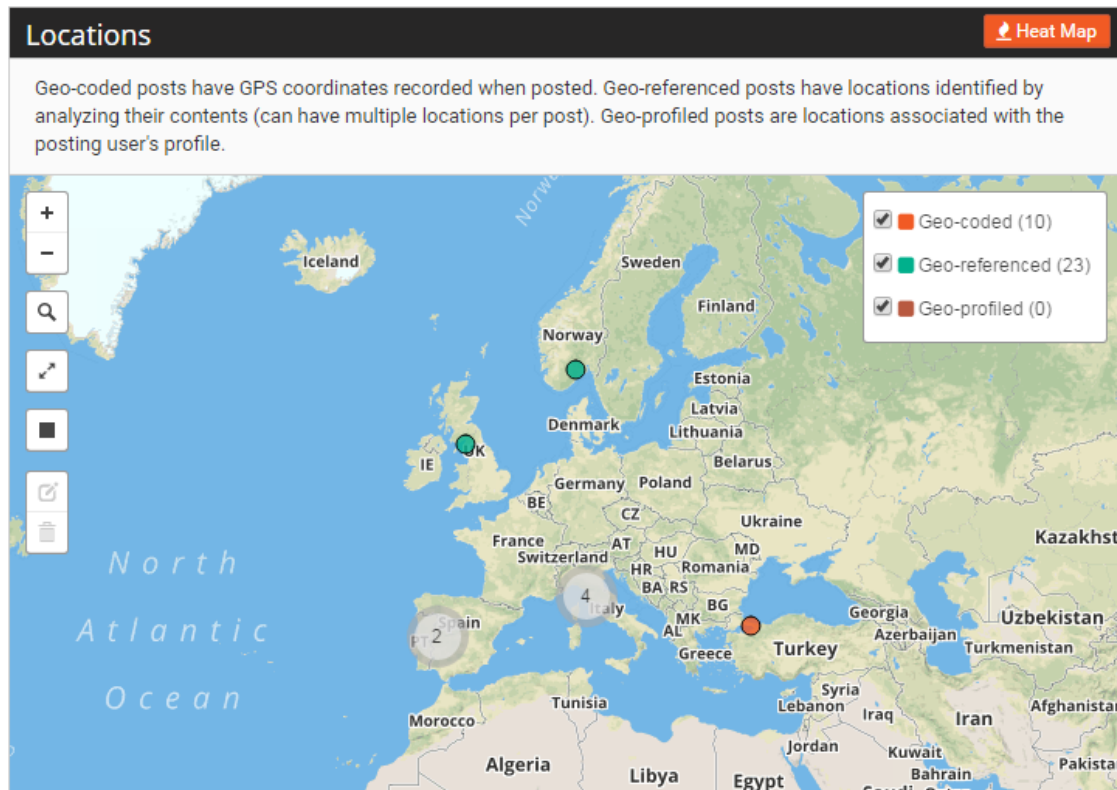
**Figure 9: Representative Map View.**

## 2.5 VK Data Collection and Analytics

### 2.5.1 VK Data Collection

VK is the largest European online social networking service. It is available in several languages, but is especially popular among Russian-speaking users. Similar to other social networks, VK allows users to message each other publicly or privately, to create groups, public pages and events, share and tag images, audio and video, and to play browser-based games.

In this reporting period, we have matured VK crawling capabilities. In particular, we have developed a capability to retrieve VK posts based on keywords. Current search grammar includes combining a set of keywords by AND'ing or OR'ing them. The search can be performed as streaming or a 1-week historical, and can be combined with a geo-location search. The search is integrated as part of Scraawl, and the UI is shown in Figure 10. Similar to other social network searches, the keywords or phrases can be entered separately, a report name can be given to the search, and either a streaming or a historical search can be chosen. In addition, a map is interactively used (not shown) to restrict the search to a geospatial region.

**Figure 10: Representative VK Search Screen.**

### 2.5.2  VK Analytics

We have also developed capabilities to compute basic statistics of the VK posts that were collected using the interface of Figure 10. In particular, top users, top words, top URLs, and top attachments along with top popular media is computed and presented to the user in one screen as shown in Figure 11. The timeline of the posts (not shown) is also presented in the same screen. Each top statistics (e.g., top words) can be further drilled down to show additional information.
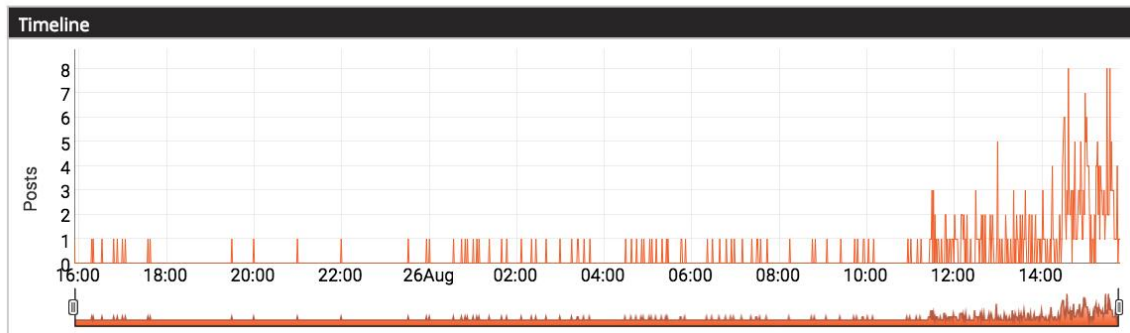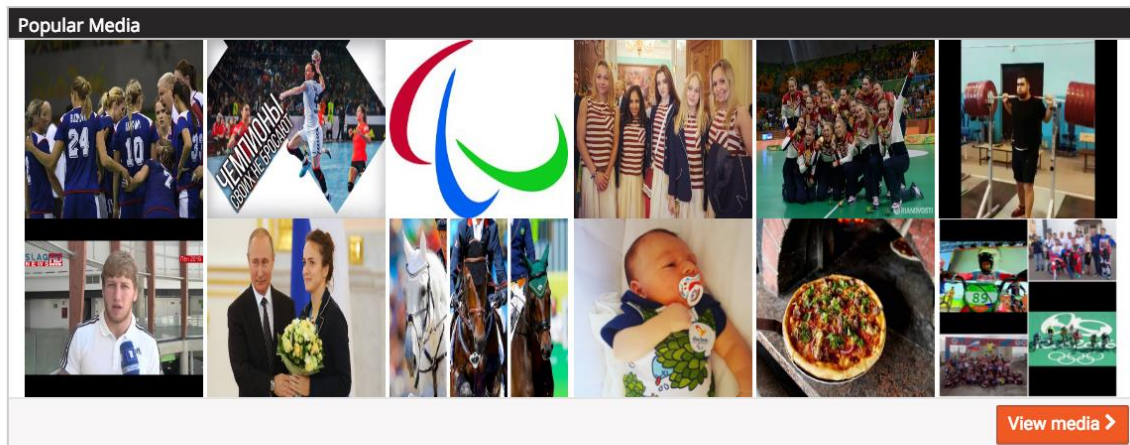
**Figure 11: Representative VK Basic Statistics View.**

Map visualization and NER have also been tailored for VK.

# 3 Documentation of Features

Please see the Scraawl help page: https://blog.scraawl.com/faq/

# 4 Conclusion

The proposed work is expected to support a variety of Navy Operational needs in the areas of Disaster Relief Operations (DRO) and assisting decision makers and field personnel to better understand and address unexpected crises. By using graph analytics at large-scales, users will be able to find key actors, extract relationships between these key actors, find pathways between leaders and followers, and interactively visualize social graphs and friendship graphs. In addition, with the added YouTube and VK capabilities, user will be able to get a broader range of textual and relationship data, and will be able to support HSCB modeling, decision making, and COA development.